

A New Distributed Transaction Protocol and Its Formal Analysis in Maude

Si Liu¹, Peter Csaba Ölveczky^{2,1}, Keshav Santhanam¹, Qi Wang¹,
Indranil Gupta¹, and José Meseguer¹

¹ University of Illinois at Urbana-Champaign

² University of Oslo

Abstract. Designers of distributed database systems face the choice between stronger consistency guarantees and better performance. A number of applications only require *read atomicity* (RA) and *prevention of lost updates* (PLU). Existing distributed database systems that meet these requirements also provide additional stronger consistency guarantees (such as *causal consistency*), and therefore incur lower performance. In this paper we define a new distributed transaction protocol, ROLA, that targets application scenarios where only RA and PLU are needed. We formally model ROLA in Maude. We then perform model checking to analyze both the correctness and the performance of ROLA. For *correctness*, we use standard model checking to analyze ROLA's satisfaction of RA and PLU. To analyze *performance* we: (a) perform statistical model checking to analyze key performance properties; and (b) compare these performance results with those obtained by also modeling and analyzing in Maude the well-known protocol Walter. Our statistical model checking results show that ROLA outperforms Walter.

1 Introduction

Distributed transaction protocols are complex distributed systems whose design is quite challenging because: (i) as for other distributed systems, validating correctness is very hard to achieve by testing alone; (ii) the high performance requirements needed in many applications are hard to measure before implementation, and expensive to compare across different implementations; and (iii) there is an unavoidable tension between the *degree of consistency* needed for the intended applications and the *high performance* required of the transaction protocol for such applications: balancing well these two requirements is essential.

In this work, we present our results on how to use formal modeling and analysis as early as possible in the design process to arrive at a mature design of a *new* distributed transaction protocol, called ROLA, meeting specific correctness and performance requirements *before* such a protocol is implemented. In this way, the above-mentioned design challenges (i)–(iii) can be adequately met. We also show how using this formal design approach it is relatively easy to *compare* ROLA with other existing transaction protocols. This is also part of meeting

design challenge (iii), since the key comparisons focus on how well each protocol balances the consistency vs. performance trade-offs for the intended applications.

ROLA in a Nutshell. Different applications require negotiating the consistency vs. performance trade-offs in different ways. The key issue is the application’s required *degree of consistency*, and how to meet such requirements with *high performance*. Cerone *et al.* [7] survey a *hierarchy of consistency models* for distributed transaction protocols including (in increasing order of strength):

- *read atomicity* (RA): either *all* or *none* of a distributed transaction’s updates are visible to another transaction (that is, there are no “fractured reads”);
- *causal consistency* (CC): if transaction T_2 is *causally dependent* on transaction T_1 , then if another transaction sees the updates by T_2 , it must also see the updates of T_1 (e.g., if A posts something on a social media, and C sees B ’s comment on A ’s post, then C must also see A ’s original post);
- *parallel snapshot isolation* (PSI): like CC but without lost updates;
- and so on, all the way up to the well-known *serializability* guarantees.

A key property of transaction protocols is the *prevention of lost updates* (PLU). The weakest consistency model in [7] satisfying both RA and PLU is PSI. However, PSI, and the well-known protocol Walter [26] implementing PSI, also guarantee CC. Cerone *et al.* conjecture that a system guaranteeing RA and PLU *without* guaranteeing CC should be useful, but up to now we are not aware of any such protocol. The point of ROLA is exactly to fill this gap: guaranteeing RA and PLU, but not CC. Two key questions that the ROLA design should answer are: (a) are there *natural applications* needing high performance where RA plus PLU provide a sufficient degree of consistency? and (b) can a new design meeting RA plus PLU *outperform* existing designs, like Walter, meeting PSI?

Regarding question (a), an example of a transaction that requires RA and PLU but not CC is the “becoming friends” transaction on social media. Bailis *et al.* [5] point out that RA is crucial for this operation: If Edinson and Neymar become friends, then Unai should not see a *fractured read* where Edinson is a friend of Neymar, but Neymar is not a friend of Edinson. An implementation of “becoming friends” must obviously guarantee PLU: the new friendship between Edinson and Neymar should not be lost. Finally, CC could be sacrificed for the sake of performance: Assume that Dani is a friend of Neymar. When Edinson becomes Neymar’s friend, he sees that Dani is Neymar’s friend, and therefore also becomes friend with Dani. The second friendship therefore causally depends on the first one. However, it does not seem crucial that others are aware of this causality: If Unai sees that Edinson and Dani are friends, then it is not necessary that he knows that (this happened *because*) Edinson and Neymar are friends.

Regarding question (b), Section 6 shows that ROLA clearly outperforms Walter in all performance requirements for all read/write transaction rates.

Maude-Based Formal Modeling and Analysis. In rewriting logic [22], distributed systems are specified as *rewrite theories*. Maude [8] is a high-performance language implementing rewriting logic and supporting various model checking

analyses. To model time and performance issues, ROLA is specified in Maude as a *probabilistic rewrite theory* [3,8]. ROLA’s RA and PLU requirements are then analyzed by standard model checking, where we disregard time issues. To estimate ROLA’s performance, and to compare it with that of Walter, we have also specified Walter in Maude, and subject the Maude models of both ROLA and Walter to *statistical model checking* analysis using the PVESTA [4] tool.

Main Contributions include: (1) the design, formal modeling, and model checking analysis of ROLA, a new transaction protocol having useful applications and meeting RA and PLU consistency properties with competitive performance; (2) a detailed performance comparison by statistical model checking between ROLA and the Walter protocol showing that ROLA outperforms Walter in all such comparisons, including higher throughput and lower average latency; (3) to the best of our knowledge the first demonstration that, by a suitable use of formal methods, a completely new distributed transaction protocol can be designed and thoroughly analyzed, as well as be compared with other designs, very early on, *before* its implementation.

2 Preliminaries

Read-Atomic Multi-Partition (RAMP) Transactions. To deal with ever-increasing amounts of data, large cloud systems *partition* their data across multiple data centers. However, guaranteeing strong consistency properties for multi-partition transactions leads to high latency. Therefore, trade-offs that combine efficiency with weaker transactional guarantees for such transactions are needed.

In [5], Bailis *et al.* propose an isolation model, *read atomic* isolation, and *Read Atomic Multi-Partition* (RAMP) transactions, that together provide efficient multi-partition operations that guarantee read atomicity (RA).

RAMP transactions use metadata and multi-versioning. Metadata is attached to each write, and the reads use this metadata to get the correct version. There are three versions of RAMP; in this paper we build on RAMP-Fast. To guarantee that all partitions perform a transaction successfully or that none do, RAMP performs two-phase writes using the two-phase commit protocol (2PC). In the *prepare* phase, each timestamped write is sent to its partition, which adds the write to its local database.³ In the *commit* phase, each such partition updates an index which contains the highest-timestamped committed version of each item stored at the partition.

RAMP assumes that there is no data *replication*: a data item is only stored at one partition. The timestamps generated by a partition P are unique identifiers but are only sequentially increasing with respect to P . A partition has access to methods $\text{GET_ALL}(I : \text{set of items})$ and $\text{PUT_ALL}(W : \text{set of } \langle \text{item, value} \rangle \text{ pairs})$.

PUT_ALL uses two-phase commit for each w in W . The first phase initiates a *prepare* operation on the partition storing $w.\text{item}$, and the second phase com-

³ RAMP does not consider write-write conflicts, so that writes are always prepared successfully (which is why RAMP does not prevent lost updates).

pletes the commit if each write partition agrees to commit. In the first phase, the client (i.e., the partition executing the transaction) passes a *version* $v : \langle \text{item}, \text{value}, ts_v, md \rangle$ to the partition, where ts_v is a timestamp generated for the transaction and md is metadata containing all other items modified in the same transaction. Upon receiving this version v , the partition adds it to a set *versions*.

When a client initiates a GET_ALL operation, then for each $i \in I$ the client will first request the latest version vector stored on the server for i . It will then look at the metadata in the version vector returned by the server, iterating over each item in the metadata set. If it finds an item in the metadata that has a later timestamp than the ts_v in the returned vector, this means the value for i is out of date. The client can then request the RA-consistent version of i .

The pseudo-code of RAMP-Fast in [5] is shown in Appendix A.

2.1 Consistency Models for Distributed Transactions

There exist a plethora of consistency models for distributed transaction systems. In [7], Cerone *et al.* characterize a number of them, including (in increasing order of strength): *read atomicity* (RA), *causal consistency* (CC: reads are consistent with transaction order; e.g., if A posts something on a social media site, and C sees B 's comment on A 's post, then C must also see A 's original post), *parallel snapshot isolation* (PSI: like CC but without lost updates), and so on.

All consistency models in [7] that satisfy RA and prevent lost updates (we will write LU for “preventing lost updates”) also satisfy CC. However, Cerone *et al.* write that “*existing consistency models do not include a counterpart of Read Atomic obtained by adding the NOCONFLICT axiom [preventing lost updates]. Such an ‘Update Atomic’ consistency model would prevent lost update anomalies without having to enforce causal consistency [...]. Update Atomic could be particularly useful when [...]*” There was until now no database design supporting such “update atomicity” (without also providing CC). Filling this gap; that is, presenting a design that does exactly that for multi-partition transactions, is what we do in this paper. (In the above hierarchy of consistency models, “update atomic” would be strictly stronger than RA, incomparable with CC, and strictly weaker than PSI.)

Walter. In [7], PSI is the weakest consistency model satisfying both RA and LU. The prototypical design supporting PSI is the *Walter* geo-replicated transactional key-value store [26].

Rewriting Logic and Maude. In rewriting logic [22] a concurrent system is specified as a *rewrite theory* $(\Sigma, E \cup A, R)$, where $(\Sigma, E \cup A)$ is a *membership equational logic theory* [8], with Σ an algebraic signature declaring sorts, subsorts, and function symbols, E a set of conditional equations, and A a set of equational axioms. It specifies the system’s state space as an algebraic data type. R is a set of *labeled conditional rewrite rules*, specifying the system’s local transitions, of the form $[l] : t \longrightarrow t' \text{ if } cond$, where $cond$ is a condition and l is a label. Such a rule specifies a transition from an instance of t to the corresponding instance of t' , provided the condition holds.

Maude [8] is a language and tool for specifying, simulating, and model checking rewrite theories. The distributed state of an object-oriented system is formalized as a *multiset* of objects and messages. A class C with attributes att_1 to att_n of sorts s_1 to s_n is declared `class C | att1 : s1, ..., attn : sn`. An object of class C is modeled as a term $\langle o : C \mid att_1 : v_1, \dots, att_n : v_n \rangle$, with o its object identifier, and where the attributes att_1 to att_n have the current values v_1 to v_n , respectively. Upon receiving a message, an object can change its state and/or send messages to other objects. For example, the rewrite rule

$$\begin{aligned} \text{rl [l] : } & \text{m}(0, z) \quad \langle 0 : C \mid a1 : x, a2 : 0' \rangle \\ \Rightarrow & \quad \langle 0 : C \mid a1 : x + z, a2 : 0' \rangle \quad \text{m}'(0', x + z) . \end{aligned}$$

defines a transition where an incoming message m , with parameters 0 and z , is consumed by the target object 0 of class C , the attribute $a1$ is updated to $x + z$, and an outgoing message $\text{m}'(0', x + z)$ is generated.

Statistical Model Checking and PVESTA. Probabilistic distributed systems can be modeled as *probabilistic rewrite theories* [3] with rules of the form

$$[l] : t(\vec{x}) \longrightarrow t'(\vec{x}, \vec{y}) \text{ if } \text{cond}(\vec{x}) \text{ with probability } \vec{y} := \pi(\vec{x})$$

where the term t' has new variables \vec{y} disjoint from the variables \vec{x} in the term t . The concrete values of the new variables \vec{y} in $t'(\vec{x}, \vec{y})$ are chosen probabilistically according to the probability distribution $\pi(\vec{x})$.

Statistical model checking [24,28] is an attractive formal approach to analyzing (purely) probabilistic systems. Instead of offering a yes/no answer, it can verify a property up to a user-specified level of confidence by running Monte-Carlo simulations of the system model. We then use PVESTA [4], a parallelization of the tool VESTA [25], to statistically model check purely probabilistic systems against properties expressed as QUATEX expressions [3]. The expected value of a QUATEX expression is iteratively evaluated w.r.t. two parameters α and δ by sampling, until we obtain a value v so that with $(1 - \alpha)100\%$ statistical confidence, the expected value is in the interval $[v - \frac{\delta}{2}, v + \frac{\delta}{2}]$.

3 The ROLA Multi-Partition Transaction Algorithm

Our new algorithm for distributed multi-partition transactions, ROLA, extends RAMP-Fast. RAMP-Fast guarantees RA, but it does not guarantee PLU since it allows a write to overwrite conflicting writes: When a partition commits a write, it only compares the write's timestamp t_1 with the local latest-committed timestamp t_2 , and updates the latest-committed timestamp with t_1 or t_2 . If the two timestamps are from two conflicting writes, then one of the writes is lost.

ROLA's key idea to prevent lost updates is to sequentially order writes on the same key from a partition's perspective by adding to each partition a data structure which maps each incoming version to an incremental sequence number. For write-only transactions the mapping can always be built; for a read-write

transaction the mapping can only be built if there has not been a mapping built since the transaction fetched the value. This can be checked by comparing the last prepared version's timestamp's mapping on the partition with the fetched version's timestamp's mapping. In this way, ROLA prevents lost updates by allowing versions to be prepared only if no conflicting prepares occur concurrently.

Algorithm 1 ROLA

Server-side Data Structures

- 1: *versions*: list of versions (item, value, timestamp ts_v , metadata md)
- 2: *latestCommit*[i]: last committed timestamp for item i
- 3: *seq*[ts]: local sequence number mapped to timestamp ts
- 4: *sqn*: local sequence counter

Server-side Methods

GET same as in RAMP-Fast

- 5: **procedure** PREPARE_UPDATE(v : version, ts_{prev} : timestamp)
- 6: $latest \leftarrow \text{last } w \in \text{versions} : w.item = v.item$
- 7: **if** $latest = \text{NULL}$ **or** $ts_{prev} = latest.ts_v$ **then**
- 8: $sqn \leftarrow sqn + 1$; $seq[v.ts_v] \leftarrow sqn$; $versions.add(v)$
- 9: **return** ACK
- 10: **else return** $latest$
- 11: **procedure** PREPARE(v : version)
- 12: $sqn \leftarrow sqn + 1$; $seq[v.ts_v] \leftarrow sqn$; $versions.add(v)$
- 13: **procedure** COMMIT(ts_c : timestamp)
- 14: $I_{ts} \leftarrow \{w.item \mid w \in \text{versions} \wedge w.ts_v = ts_c\}$
- 15: **for** $i \in I_{ts}$ **do**
- 16: **if** $seq[ts_c] > seq[latestCommit[i]]$ **then** $latestCommit[i] \leftarrow ts_c$

Coordinator-side Methods

PUT_ALL, GET_ALL same as in RAMP-Fast

- 17: **procedure** UPDATE(I : set of items, OP : set of operations)
 - 18: $ret \leftarrow \text{GET_ALL}(I)$; $ts_{tx} \leftarrow \text{generate new timestamp}$
 - 19: **parallel-for** $i \in I$ **do**
 - 20: $ts_{prev} \leftarrow ret[i].ts_v$; $v \leftarrow ret[i].value$
 - 21: $w \leftarrow \langle item = i, value = op_i(v), ts_v = ts_{tx}, md = (I - \{i\}) \rangle$
 - 22: $p \leftarrow \text{PREPARE_UPDATE}(w, ts_{prev})$
 - 23: **if** $p = latest$ **then**
 - 24: invoke application logic to, e.g., abort and/or retry the transaction
 - 25: **end parallel-for**
 - 26: **parallel-for** server s : s contains an item in I **do**
 - 27: invoke COMMIT(ts_{tx}) on s
 - 28: **end parallel-for**
-

More specifically, ROLA adds two partition-side data structures: sqn , denoting the local sequence counter, and $seq[ts]$, that maps a timestamp to a local sequence number. ROLA also changes the data structure of $versions$ in RAMP from a set to a list. ROLA then adds two methods: the coordinator-side⁴ method $UPDATE(I : \text{set of items}, OP : \text{set of operations})$ and the partition-side method $PREPARE_UPDATE(v : \text{version}, ts_{prev} : \text{timestamp})$ for read-write transactions. Furthermore, ROLA changes two partition-side methods in RAMP: $PREPARE$, besides adding the version to the local store, maps its timestamp to the increased local sequence number; and $COMMIT$ marks versions as committed and updates an index containing the highest-sequenced-timestamped committed version of each item. These two partition-side methods apply to both write-only and read-write transactions. ROLA invokes RAMP-Fast's PUT_ALL , GET_ALL and GET methods to deal with read-only and write-only transactions.

ROLA starts a read-write transaction with the $UPDATE$ procedure. It invokes RAMP-Fast's GET_ALL method to retrieve the values of the items the client wants to update, as well as their corresponding timestamps. ROLA writes then proceed in two phases: a first round of communication places each timestamped write on its respective partition. The timestamp of each version obtained previously from the GET_ALL call is also packaged in this *prepare* message. A second round of communication marks versions as committed.

At the partition-side, the partition begins the $PREPARE_UPDATE$ routine by retrieving the last version in its $versions$ list with the same item as the received version. If such a version is not found, or if the version's timestamp ts_v matches the passed-in timestamp ts_{prev} , then the version is deemed prepared. The partition keeps a record of this locally by incrementing a local sequence counter and mapping the received version's timestamp ts_v to the current value of the sequence counter. Finally the partition returns an ACK to the client. If ts_{prev} does not match the timestamp of the last version in $versions$ with the same item, then this *latest* timestamp is simply returned to the coordinator.

If the coordinator receives an ACK from $PREPARE_UPDATE$, it immediately commits the version with the generated timestamp ts_{tx} . If the returned value is instead a timestamp, the transaction is aborted.

Example. Consider the ROLA execution depicted in Figure 1, where two read-write transactions $T_1 : r(y); w(x_1); w(y_1)$ and $T_2 : r(y); w(y_2)$ are attempting concurrent writes, and a read-only transaction $T_3 : r(x); r(y)$ proceeds while T_1 is writing. T_1 and T_2 read the same version y_0 . Both T_1 and T_2 perform the two-phase commit protocol on two partitions, P_x and P_y . However, T_2 fails to prepare y_2 after T_1 has prepared y_1 , because when T_2 's prepare arrives at P_y , the timestamp of the last version store on P_y is 1, which is not equal to $ts_{prev} = 0$ in T_2 's prepare. T_2 , upon receiving the returned version y_1 , could abort the transaction or retry with a new transaction on y_1 . Either way the lost update problem is avoided. Regarding the case with T_1 and T_3 , T_3 reads from P_x after P_x has committed T_1 's write to x , but T_3 reads from P_y before P_y has committed

⁴ The *coordinator*, or *client*, is the partition executing the transaction.

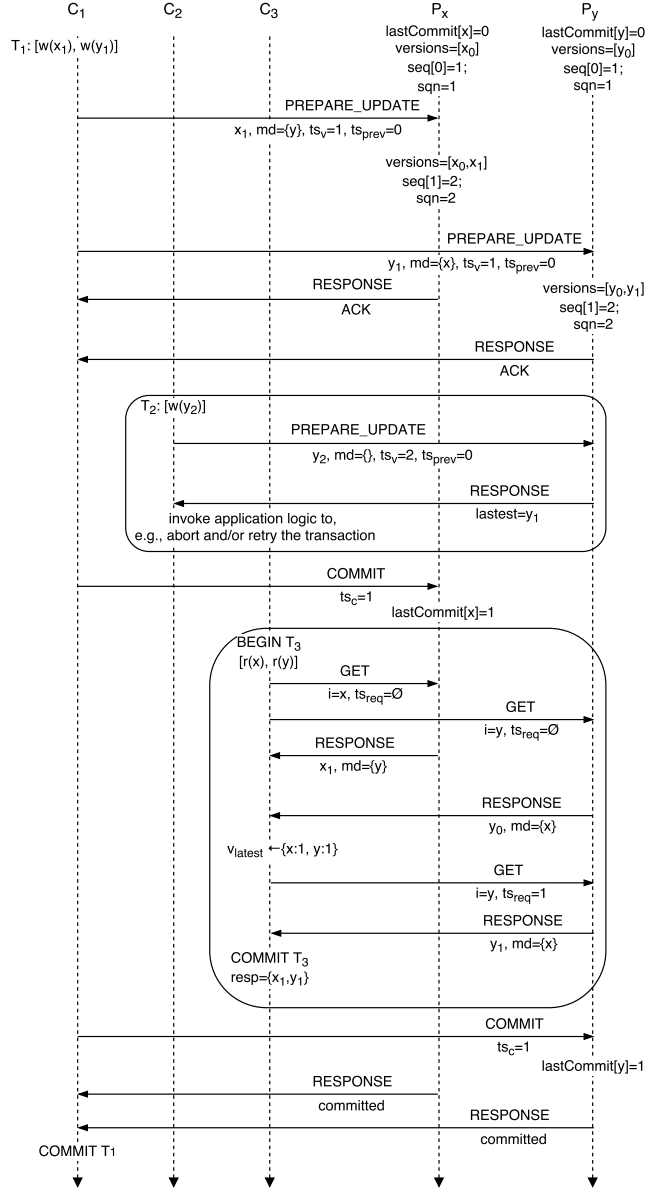


Fig. 1. Example execution of ROLA with three transactions

T_1 's write to y . Thus, T_3 's first-round reads would violate RA if it returns them. Using the metadata attached to its first-round reads, T_3 determines to issue a second-round read to fetch the missing data from P_y . After completing the

second-round read, T_3 can safely return T_1 's writes, not violating RA. Note that in this example RAMP allows T_2 to commit, thus overwriting T_1 's writes, which are then lost.

Why it works. ROLA uses a two-phase commit protocol in order to detect concurrent writes. The first phase declares an intent to commit a write at the partition. The partition can accept a preparation if there is no other prepared version after the latest commit associated with the incoming preparation. This in effect imposes a total order on the preparations, and thus on the commits, from the partition's perspective. In other words, the partition sees no logically concurrent updates. Our algorithm therefore provides read atomicity, and prevents updates from being lost, as concurrent updates are a necessary condition for lost updates.

By leveraging the partition-side sequence counter to commit, ROLA not only prevents lost updates, but also make writes progress at the partition-side, and thus more recent prepared version can be reflected (we refer to this as ROLA's *progressing property*). This is different from RAMP-F where later prepared writes may never be fetched by reads as *latestCommit* only updates by simply comparing the coordinator-side timestamps.

4 A Probabilistic Model of ROLA

This section defines a formal executable probabilistic model of ROLA. The whole specification is given at <https://sites.google.com/site/fase18submission/>.

As mentioned in Section 2, statistical model checking assumes that the system is *fully probabilistic*; that is, has no unquantified nondeterminism. We follow the techniques in [9] to obtain such a model. The key idea is that message delays are sampled probabilistically from dense/continuous time intervals. The probability that two messages will have the same delay is therefore 0. If events only take place when a message arrives, then two events will not happen at the same time, and therefore unquantified nondeterminism is eliminated.

We are also interested in correctness analysis of a model that captures all possible behaviors from a given initial configuration. We obtain such a non-deterministic untimed model, that can be subjected to standard model checking analysis, by just removing all message delays from our probabilistic timed model.

4.1 Probabilistic Sampling

Nodes send messages of the form $[\Delta, rcvr \leftarrow msg]$, where Δ is the message delay, $rcvr$ is the recipient, and msg is the message content. When time Δ has elapsed, this message becomes a *ripe* message $\{T, rcvr \leftarrow msg\}$, where T is the "current global time" (used for analysis purposes only). Such a ripe message must be consumed by the receiver $rcvr$ before time advances. Δ can be sampled from certain distributions (lognormal, Weibull, etc.) for statistical model checking, or, as mentioned above, removed (or set to zero) for correctness analysis.

To sample message delays from different distributions, we use the following functionality provided by Maude: The built-in function `random`, where `random(k)` returns the k -th pseudo-random number as a number between 0 and $2^{32} - 1$, and the built-in constant `counter` with an (implicit) rewrite rule `counter => N:Nat`. The first time `counter` is rewritten, it rewrites to 0, the next time it rewrites to 1, and so on. Therefore, each time `random(counter)` rewrites, it rewrites to the next random number. Since Maude does not rewrite `counter` when it appears in the condition of a rewrite rule, we encode a probabilistic rewrite rule $t(\vec{x}) \longrightarrow t'(\vec{x}, \vec{y})$ if *cond*(\vec{x}) with probability $\vec{y} := \pi(\vec{x})$ in Maude as the rule $t(\vec{x}) \longrightarrow t'(\vec{x}, \text{sample}(\pi(\vec{x})))$ if *cond*(\vec{x}). The following operator `sampleLogNormal` is used to sample a value from a lognormal distribution with mean `MEAN` and standard deviation `SD`:

```
op sampleLogNormal : Float Float -> [Float] .
eq sampleLogNormal(MEAN,SD) = exp(MEAN + SD * sampleNormal) .

op sampleNormal : -> [Float] .    op sampleNormal : Float -> [Float] .
eq sampleNormal = sampleNormal(float(random(counter) / 4294967296)) .
eq sampleNormal(RAND) = sqrt(- 2.0 * log(RAND)) * cos(2.0 * pi * RAND) .
```

`random(counter) / 4294967296` rewrites to a different “random” number between 0 and 1 each time it is rewritten, and this is used to define the sampling function. For example, the message delay `rd` to a remote site can be sampled from a lognormal distribution with mean 3 and standard deviation 2 as follows:

```
eq rd = sampleLogNormal(3.0, 2.0) .
```

4.2 Data Types, Classes, and Messages

We formalize ROLA in an object-oriented style, where the state consists of a number of *partition* objects, each modeling a partition of the database, and a number of messages traveling between the objects. A *transaction* is formalized as an object which resides inside the partition object that executes the transaction.

Data Types. A *version* is a timestamped version of a data item (or key) and is modeled as a 4-tuple `version(key,value,timestamp,metadata)`. consisting of the key, its value, and the version’s timestamp and metadata. A timestamp is modeled as a pair `ts(addr,sgn)` consisting of a partition’s identifier *addr* and a local sequence number *sgn*. that together uniquely identify a write transaction. Metadata are modeled as a set of keys, denoting, for each key, the other keys that are written in the same transaction. For example, if a transaction writes keys *x*, *y*, and *z*, then versions of *x* have as metadata the set $\{y, z\}$.

```
sorts Key Value Timestamp Version .
pr SET{Key} * (sort Set{Key} to KeySet) .

op ts : Address Nat -> Timestamp .
op version : Key Value Timestamp KeySet -> Version [ctor] .
```

Terms of sort `KeySet` are terms k_1, k_2, \dots, k_j , with each k_i a term of sort `Key`.

The sort `OperationList` represents lists of read and write operations as terms such as $(x := \text{read } k1) (y := \text{read } k2) \text{write}(k1, x+y)$, where `LocalVar` denotes the “local variable” that stores the value of the key read by the operation, and `Expression` is an expression involving the transaction’s local variables:

```
op write : Key Expression -> Operation [ctor] .
op _:=read_ : LocalVar Key -> Operation [ctor] .
pr LIST{Operation} * (sort List{Operation} to OperationList) .
```

Classes. A *transaction* is modeled as an object of the following class `Txn`:

```
class Txn | operations : OperationList, readSet : Versions,
          localVars : LocalVars,       latest : KeyTimestamps .
```

The `operations` attribute denotes the transaction’s operations. The `readSet` attribute denotes the versions read by the read operations. `localVars` maps the transaction’s local variables to their current values. `latest` stores the local view as a mapping from keys to their respective latest committed timestamps.

A *partition* (or *site*) stores parts of the database, and executes the transactions for which it is the coordinator/server. A partition is formalized as an object instance of the following class `Partition`:

```
class Partition | datastore : Versions,      sqn : Nat,
                gotTxns : ObjectList,       executing : Object,
                committed : ObjectList,     aborted : ObjectList,
                tsSqn : TimestampSqn,      latestCommit : KeyTimestamps,
                votes : Vote,               voteSites : TxnAddrSet,
                1stGetSites : TxnAddrSet, 2ndGetSites : TxnAddrSet,
                commitSites : TxnAddrSet .
```

The `datastore` attribute represents the partition’s local database as a list of versions for each key stored at the partition. The attribute `latestCommit` maps to each key the timestamp of its last committed version. `tsSqn` maps each version’s timestamp to a local sequence number `sqn`. The attributes `gotTxns`, `executing`, `committed` and `aborted` denote the transaction(s) which are, respectively, waiting to be executed, currently executing, committed, and aborted. A partition executes transactions sequentially. Concurrent transactions can be modeled by multiple transactions executed by different partitions.

The attribute `votes` stores the votes, as triples $\text{vote}(txn, part, result)$, from the partitions which participate in the two-phase commit. The remaining attributes denote the partitions from which the executing partition is awaiting votes, committed acks, first-round get replies, and second-round get replies. These are represented by terms $\text{addrs}(txn, \text{setOfPartitions})$.

The state also contains a “table” mapping each data item to the partition storing the item.

The following shows an initial state (with some parts replaced by ‘...’) with two partitions, `p1` and `p2`, that are coordinators for, respectively, transactions `t1`,

and **t2** and **t3**. **p1** stores the data items **x** and **z**, and **p2** stores **y**. Transaction **t1** is the read-only transaction (**x1 :=read x**) (**y1 :=read y**), transaction **t2** is a write-only transaction **write(y, 3)** **write(z, 8)**, while transaction **t3** is a read-write transaction on data item **x**. The states also include a buffer of messages in transit and the global clock value, and a table which assigns to each data item the site storing the item. Initially, the value of each item is **[0]**; the version's timestamp is empty (**eptTS**), and metadata is an empty set.

```
eq init = { 0.0 | nil}
< tb : Table | table : [sites(x,p1) ;; sites(y,p2) ;; sites(z,p1)] >
< p1 : Partition |
    gotTxns: < t1 : Txn | operations: ((x1 :=read x) (y1 :=read y)),
              readSet: empty, latest: empty,
              localVar: (x1 |-> [0], y1 |-> [0]) >,
    datastore: (version(x, [0], eptTS, empty)
                version(z, [0], eptTS, empty)),
    sqn: 1, ... >
< p2 : Partition |
    gotTxns: < t2 : Txn | operations: (write(y, 3) write(z, 8)), ... >
              < t3 : Txn | operations: ((x1 := read x)
                                         write(x, x1 plus 1)), ... >
    datastore: version(y, [0], eptTS, empty), ... > .
```

Messages. The message **prepare(txn, version, sender)** sends a version from a write-only transaction to its partition, and **prepare(txn, version, ts, sender)** does the same thing for other transactions, with **ts** the timestamp of the version it has read. The partition replies with a message **prepare-reply(txn, vote, sender)**, where **vote** tells whether this partition can commit the transaction. A message **commit(txn, ts, sender)** marks the versions with timestamp **ts** as committed. **get(txn, key, ts, sender)** asks for the highest-timestamped committed version or a missing version for **key** by timestamp **ts**, and **response1(txn, version, sender)** and **response2(txn, version, sender)** respond to first/second-round **get** requests.

4.3 Formalizing ROLA's Behaviors

This section formalizes the dynamic behaviors of ROLA using rewrite rules, referring to the corresponding lines in Algorithm 1.⁵

Starting a write-only transaction (lines 17 – 22). A partition starts executing a transaction by moving the first transaction (TID) in **gotTxns** to **executing**. If it is a write-only transaction (**write-only(OPS)**), the partition: (i) uses the function **genPuts** to generate all **prepare** messages; (ii) uses a function **prepareSites** to remember the sites **PIDS** from which it awaits votes for transaction TID in the **voteSites** attribute; and (iii) increments its local sequence number by one:⁶

⁵ We do not give variable declarations, but follow the convention that variables are written in (all) capital letters.

⁶ The variables **AS** and **AS'** denote the “remaining” attributes in the two objects.

```

crl [start-wo-txn] :
  < TABLE : Table | table: PARTITION-TABLE >
  < PID : Partition |
    gotTxns: (< TID : Txn | operations: OPS, localVars: VARS, AS >
              ;; TXNS),
    executing: noTxn, sqn: SQN, voteSites: VSTS, AS' >
=>
  < TABLE : Table | table: PARTITION-TABLE >
  < PID : Partition |
    gotTxns: TXNS,
    executing: < TID : Txn | operations: OPS, localVars: VARS, AS >,
    sqn: SQN', voteSites: (VSTS ; addrS(TID,PIDS)), AS' >
  genPuts(OPS,PID,TID,SQN',VARS,PARTITION-TABLE)
  if SQN' := SQN + 1 /\ write-only(OPS) /\
    PIDS := prepareSites(OPS,PID,TRANSITION-TABLE) .

```

Otherwise, if the first transaction in `gotTxns` is a read-only or read-write transaction, the replica updates `1stGetSites` instead to keep track of the replicas from which it receives the versions from the first-round gets. Similarly, the expression `genGets` generates all get messages for the keys concerned by `TID`.⁷ The expression `1stSites` gives the corresponding replicas for those keys:

```

crl [start-ro-or-rw-txn] :
  < TABLE : Table | table: REPLICIA-TABLE >
  < RID : Partition | gotTxns: (< TID : Txn | operations: OPS,
                                latest: empty, AS > ;; TXNS),
    executing: noTxn,
    1stGetSites: 1STGETS, AS' >
=>
  < TABLE : Table | table: REPLICIA-TABLE >
  < RID : Partition | gotTxns: TXNS,
    executing: < TID : Txn | operations: OPS,
                                latest: vl(OPS), AS >,
    1stGetSites: (1STGETS ; addrS(TID,RIDS)), AS' >
  genGets(OPS,RID,TID,REPLICIA-TABLE)
  if (not write-only(OPS)) /\
    RIDS := 1stSites(OPS,RID,REPLICIA-TABLE) .

```

*Receiving **prepare** messages (lines 5–10).* When a partition receives a **prepare** message for a read-write transaction, the partition first determines whether the timestamp of the last version (`VERSION`) in its local version list `VS` matches the incoming timestamp `TS'` (which is the timestamp of the version read by the transaction). If so, the incoming version is added to the local store, the map `tsSqn` is updated, and a positive reply (**true**) to the prepare message is sent (“**return ack**” in our pseudo-code); otherwise, a negative reply (**false**, or

⁷ In this paper we consider one-shot reads [14] that do not include cross-server key dependencies. Thus one-shot reads can be issued in parallel.

“**return latest**” in the pseudo-code) is sent. Depending on whether the sender PID' of the *prepare* message happens to be PID itself, the reply is equipped with a local message delay ld or a remote message delay rd , both of which are sampled probabilistically from distributions with different parameters:⁸

```

crl [receive-prepare-rw] :
  {T, PID <- prepare(TID, version(K, V, TS, MD), TS', PID')}
  < PID : Partition | datastore: VS, sqn: SQN, tsSqn: TSSQN, AS' >
=>
  if VERSION == eptVersion or tstamp(VERSION) == TS'
  then < PID : Partition | datastore: (VS version(K, V, TS, MD)), sqn: SQN',
        tsSqn: insert(TS, SQN', TSSQN), AS' >
    [if PID == PID' then ld else rd fi,
     PID' <- prepare-reply(TID, true, PID)]
  else < PID : Partition | datastore: VS, sqn: SQN, tsSqn: TSSQN, AS' >
    [if PID == PID' then ld else rd fi,
     PID' <- prepare-reply(TID, false, PID)] fi
  if SQN' := SQN + 1 /\ VERSION := latestPrepared(K, VS) .

```

In instead the received *prepare* message was for a write-only transaction, the replica simply adds the received version to its local datastore, and maps the associated timestamp to the incremented sequence number (by the function $insert(TS, SQN', TSSQN)$). Depending on whether the message sender RID' is the replica itself, the out-going message is equipped with a local message delay ld or a remote message delay rd . Both delays, as mentioned before, are sampled on a certain distribution. Note that case (i) always considers successful preparations.

```

crl [receive-prepare-wo] :
  < RID : Partition | datastore: VS,
    sqn: SQN,
    tsSqn: TSSQN, AS' >
  {T, RID <- prepare(TID, version(K, V, TS, MD), RID')}
=>
  < RID : Partition | datastore: (VS version(K, V, TS, MD)),
    sqn: SQN',
    tsSqn: insert(TS, SQN', TSSQN), AS' >
  [if RID == RID' then ld else rd fi, RID' <- prepare-reply(TID, true, RID)]
  --- always "true" for write-only prepare
  if SQN' := SQN + 1 .

```

In case (ii), the replica first determines whether the timestamp of the last $VERSION$ in the local version list VS matches the incoming timestamp TS' which is the timestamp associated with the version fetched by the previous get in the same read-write transaction. If matched, the incoming version is added to the local datastore; otherwise, a negative ack (denoted by **false**) is sent back:

⁸ The variable AS' denotes the “remaining” attributes in the object.

```

cr1 [receive-prepare-rw] :
  < RID : Partition | datastore: VS,
    sqn: SQN,
    tsSqn: TSSQN, AS' >
  {T, RID <- prepare(TID,version(K,V,TS,MD),TS',RID')}
=>
  if VERSION == eptVersion or tstamp(VERSION) == TS'
    then < RID : Partition | datastore: (VS version(K,V,TS,MD)),
      sqn: SQN',
      tsSqn: insert(TS,SQN',TSSQN), AS' >
      [if RID == RID' then ld else rd fi,
       RID' <- prepare-reply(TID,true,RID)]
    else < RID : Partition | datastore: VS,
      sqn: SQN,
      tsSqn: TSSQN, AS' >
      [if RID == RID' then ld else rd fi,
       RID' <- prepare-reply(TID,false,RID)]
  fi
  if SQN' := SQN + 1 /\
    VERSION := latestPrepared(K,VS) .

```

Receiving negative replies (lines 23–24). When a site receives a **prepare-reply** message with vote **false**, it aborts the transaction by moving it to the **aborted** list, and removes **PID'** from the “vote waiting list” for this transaction: If the transaction has been aborted, the incoming **prepare-reply** message is simply consumed by the replica whether it is a negative ack (**FLAG** is a boolean variable).

```

r1 [receive-prepare-reply-false-executing] :
  {T, PID <- prepare-reply(TID, false, PID')}
  < PID : Partition | executing: < TID : Txn | AS >, aborted: TXNS,
    voteSites: VSTS addrS(TID, (PID', PIDS)), AS' >
=>
  < PID : Partition | executing: noTxn,
    aborted: (TXNS ;; < TID : Txn | AS >),
    voteSites: VSTS addrS(TID, PIDS), AS' > .

```

```

r1 [receive-prepare-reply-aborted] :
  < RID : Partition | aborted: (TXNS ;; < TID : Txn | AS > ;; TXNS'),
    voteSites: VSTS, AS' >
  {T, RID <- prepare-reply(TID,FLAG,RID')} --- no matter what FLAG is
=>
  < RID : Partition | aborted: (TXNS ;; < TID : Txn | AS > ;; TXNS'),
    voteSites: remove(TID,RID',VSTS), AS' > .

```

Receiving Acks (Lines 26–28). Upon receiving a “true” vote, the replica first checks if all votes has now been collected. The expression **VSTS'[TID]** projects for **TID** the remaining replicas from which it is waiting for the votes. If all (“yes”)

votes received, the replica starts to commit TID on the associated replicas by invoking `genCommits` to generate all commit messages with the commit timestamp including the current sequence number `sqn`). The replica also adds to `commitSites` the replicas from which it is waiting for the committed messages to commit the transaction:

```
crl [receive-prepare-reply-true-executing] :
  < TABLE : Table | table: REPLICIA-TABLE >
  < RID : Partition | executing: < TID : Txn | operations: OPS, AS >,
    voteSites: VSTS, sqn: SQN,
    commitSites: CMTS, AS' >
  {T, RID <- prepare-reply(TID,true,RID')}
=>
  < TABLE : Table | table: REPLICIA-TABLE >
  if VSTS'[TID] == empty --- all votes received and all yes!
    then < RID : Partition | executing: < TID : Txn | operations: OPS, AS >,
      voteSites: VSTS', sqn: SQN,
      commitSites: (CMTS ; addrs(TID,RIDS)), AS' >
      genCommits(TID,SQN,RIDS,RID)
    else < RID : Partition | executing: < TID : Txn | operations: OPS, AS >,
      voteSites: VSTS', sqn: SQN,
      commitSites: CMTS, AS' >
  fi
  if VSTS' := remove(TID,RID',VSTS) /\
    RIDS := commitSites(OPS,RID,REPLICIA-TABLE) .
```

Receiving commit messages (lines 13–16). Upon receiving a commit message, the partition invokes the function `cmt` to commit the transaction. `cmt` looks up `tsSqn` for the commit timestamp `TS` and latest committed version's timestamp in `LC`, and updates the latest committed version if `TS`'s local sequence number is higher. A committed message is then sent back to confirm the commit:

```
r1 [receive-commit] :
  {T, PID <- commit(TID,TS,PID')}
  < PID : Partition | tsSqn: TSSQN, dataStore: VS, latestCommit: LC, AS' >
=>
  < PID : Partition | tsSqn: TSSQN, dataStore: VS,
    latestCommit: cmt(LC,VS,TSSQN,TS), AS' >
  [if PID == PID' then ld else rd fi, PID' <- committed(TID,PID)] .
```

Receiving Committed Messages. Upon receiving a committed message, the replica first checks if all committed messages has now been collected. The expression `CMTS'[TID]` projects for `TID` the remaining replicas from which it is waiting for the committed messages. If the projection is `empty`, the replica commits the transaction:

```
crl [receive-committed] :
```



```

< RID : Partition | executing: < TID : Txn | AS >,
    committed: TXNS, commitSites: CMTS, AS' >
{T, RID <- committed(TID,RID') }
=>
if CMTS'[TID] == empty --- all "committed" received
  then < RID : Partition | executing: noTxn,
      committed: (TXNS ;; < TID : Txn | AS >),
      commitSites: CMTS', AS' >
  else < RID : Partition | executing: < TID : Txn | AS >,
      committed: TXNS,
      commitSites: CMTS', AS' >
fi
if CMTS' := remove(TID,RID',CMTS) .

```

Receiving Get Messages. Upon receiving a `get` message, depending of the associated timestamp TS (if TS is an empty timestamp `eptTS`, the incoming message is the first-round `get`; otherwise, it is the second-round `get`), the replica replies with the corresponding version determined by the function `vmatch`. For the first-round `get`, `vmatch` looks up LC for the latest committed version; for the second-round `get`, `vmatch` returns the matched timestamped version of TS:

```

r1 [receive-get] :
  < RID : Partition | datastore: VS, latestCommit: LC, AS' >
  {T, RID <- get(TID,K,TS,RID') }
=>
  < RID : Partition | datastore: VS, latestCommit: LC, AS' >
  [if RID == RID' then ld else rd fi,
   RID' <- (if TS == eptTS then response1(TID,vmatch(K,VS,LC),RID)
   else response2(TID,vmatch(K,VS,TS),RID) fi)] .

```

Receiving Replies to First-Round Gets. Upon receiving a returned version for the first-round `get`, the replica adds it to the read set, and updates `localVars` correspondingly. Once the replica has collected all replies to the first-round gets, it then determines whether a second-round `get` is needed. The expression `gen2ndGets(TID,VL',RS',RID,REPLICA-TABLE)` generates possible second-round `get` messages based on the updated `latest`, `VL'`, and `readSet`, `RS'`. In the case that a second-round `get` is not needed, `gen2ndGets` generates no messages, and `RIDS` is an empty set. Note that `RS'` is specially needed if `TID` is a read-write transaction:

```

cr1 [receive-response1] :
  < TABLE : Table | table: REPLICA-TABLE >
  < RID : Partition | executing: < TID : Txn |
      operations: (OPS (X :=read K) OPS'),
      readSet: RS, localVars: VARS, latest: VL, AS >,
      1stGetSites: 1STGETS,
      2ndGetSites: 2NDGETS, AS' >

```

```

    {T, RID <- response1(TID,version(K,V,TS,MD),RID')}
=>
    < TABLE : Table | table: REPLICAS-TABLE >
    if 1STGETS'[TID] == empty
    then < RID : Partition | executing: < TID : Txn |
        operations: (OPS (X :=read K) OPS'),
        readSet: RS', localVars: insert(X,V,VARS),
        latest: VL', AS >,
        1stGetSites: 1STGETS',
        2ndGetSites: (2NDGETS ; addrs(TID,RIDS)), AS' >
        gen2ndGets(TID,VL',RS',RID,REPLICAS-TABLE)
    else < RID : Partition | executing: < TID : Txn |
        operations: (OPS (X :=read K) OPS'),
        readSet: RS', localVars: insert(X,V,VARS),
        latest: VL', AS >,
        1stGetSites: 1STGETS',
        2ndGetSites: 2NDGETS, AS' >

    fi
    if RS' := RS version(K,V,TS,MD) /\
        VL' := lat(VL,MD,TS) /\
        1STGETS' := remove(TID,RID',1STGETS) /\
        RIDS := 2ndSites(VL',RS',RID,REPLICAS-TABLE) .

```

Receiving Replies to Second-Round Gets. Upon receiving a returned version for the second-round get, the replica simply overwrites the version fetched by the first-round get. It then updates the local variables and the remaining replicas from which it is waiting for the second-round gets:

```

r1 [receive-response2] :
    < RID : Partition | executing: < TID : Txn |
        operations: (OPS (X :=read K) OPS'),
        readSet: (RS version(K,V',TS',MD') RS'),
        localVars: VARS, AS >,
        2ndGetSites: 2NDGETS, AS' >
    {T, RID <- response2(TID,version(K,V,TS,MD),RID')}
=>
    < RID : Partition | executing: < TID : Txn |
        operations: (OPS (X :=read K) OPS'),
        readSet: (RS version(K,V,TS,MD) RS'),
        localVars: insert(X,V,VARS), AS >,
        2ndGetSites: remove(TID,RID',2NDGETS), AS' > .

```

Committing Reads. Once the replica has no remaining replicas from which it is waiting for the replies to either first-round gets or second-round gets (denoted by `1STGETS[TID] == empty` and `2NDGETS[TID] == empty`), it starts to commit the reads. There are two cases to consider: (i) a read-only transaction; or (ii) a read-write transaction.

In case (i), the replica simply puts TID in `committed`; in case (ii), the replica further generates all prepare messages for each version concerned with newly generated timestamp including the incremented sequence number `SQN'`. The prepared versions are computed based on the previously fetched reads reflected in `VARs`, and the prepare messages also include the timestamps of the previously fetched reads in `RS`:

```

crl [commit-reads] :
  < TABLE : Table | table: REPLICIA-TABLE >
  < RID : Partition | executing: < TID : Txn | operations: OPS,
                                localVars: VARs,
                                readSet: RS, AS >,
                                committed: TXNS, 1stGetSites: 1STGETS,
                                2ndGetSites: 2NDGETS, sqn: SQN, voteSites: VSTS, AS' >
=>
  < TABLE : Table | table: REPLICIA-TABLE >
  if read-only(OPS)
    then < RID : Partition | executing: noTxn,
          committed: (TXNS ;; < TID : Txn |
                      operations: OPS, localVars: VARs,
                      readSet: RS, AS >),
          1stGetSites: 1STGETS, 2ndGetSites: 2NDGETS,
          sqn: SQN, voteSites: VSTS, AS' >
    else < RID : Partition | executing: < TID : Txn | operations: OPS,
          localVars: VARs, readSet: RS, AS >,
          committed: TXNS, 1stGetSites: 1STGETS,
          2ndGetSites: 2NDGETS, sqn: SQN',
          voteSites: (VSTS ; voteSites(TID,RIDS)), AS' >
          genPuts(OPS,RID,TID,SQN',VARs,RS,REPLICIA-TABLE)
  fi
  if 1STGETS[TID] == empty /\ 2NDGETS[TID] == empty /\
    SQN' := SQN + 1 /\ RIDS := prepareSites(OPS,RID,REPLICIA-TABLE) .

```

5 Correctness Analysis of ROLA

In this section we use reachability analysis to analyze whether ROLA guarantees read atomicity and prevents lost updates.

For both correctness and performance analysis, we add to the state an object

```
< m : Monitor | log: log >
```

which stores crucial information about each transaction. The *log* is a list of records `record(tid, issueTime, finishTime, reads, writes, committed)`, with *tid* the transaction's ID, *issueTime* its issue time, *finishTime* its commit/abort time, *reads* the versions read, *writes* the versions written, and *committed* a flag that is `true` if the transaction is committed.

We modify our model by updating the **Monitor** when needed. For example, when the coordinator has received all **committed** messages, the monitor records the commit time (T) for that transaction, and sets the “committed” flag to **true**⁹:

```

cr1 [receive-committed] :
  {T, PID <- committed(TID, PID')}
  < M : Monitor | log: (LOG record(TID, T', T'', RS, WS, false) LOG') >
  < PID : Partition | executing: < TID : Txn | AS >,
    committed: TXNS, commitSites: CMTS, AS' >
=>
  if CMTS'[TID] == empty --- all "committed" received
  then < M : Monitor | log: (LOG record(TID, T', T, RS, WS, true) LOG') >
    < PID : Partition | executing: noTxn, commitSites: CMTS',
      committed: (TXNS ;; < TID : Txn | AS >, AS' >
  else < M : Monitor | log: (LOG record(TID, T', T'', RS, WS, false) LOG') >
    < PID : Partition | executing: < TID : Txn | AS >,
      committed: TXNS, commitSites: CMTS', AS' > fi
  if CMTS' := remove(TID, PID', CMTS) .

```

Since ROLA is terminating if a finite number of transactions are issued, we analyze the different (correctness and performance) properties by inspecting this monitor object in the final states, when all transactions are finished.

Read Atomicity. A system guarantees RA if it prevents fractured reads, and also prevents transactions from reading uncommitted, aborted, or intermediate data [5], where a transaction T_j exhibits *fractured reads* if transaction T_i writes version x_m and y_n , T_j reads version x_m and version y_k , and $k < n$ [5].

We analyze this property by searching for a reachable *final* state (arrow $\Rightarrow!$) where the property does *not* hold:

```

search [1] initConfig =>! C:Config < M:Address : Monitor | log: LOG:Record >
  such that fracRead(LOG) or abortedRead(LOG) .

```

The function **fracRead** checks whether there are fractured reads in the execution log. There is a fractured read if a transaction TID2 reads X and Y, transaction TID1 writes X and Y, TID2 reads the version TSX of X written by TID1, and reads a version TSY' of Y written *before* TSY (TSY' < TSY). Since the transactions in the log are ordered according to start time, TID2 could appear *before* or *after* TID1 in the log. We spell out the case when TID1 comes before TID2:

```

op fracRead : Record -> Bool .
ceq fracRead(LOG ;
  record(TID1, T1, T1', RS1, (version(X, VX, TSX, MDX), version(Y, VY, TSY, MDY)), true) ; LOG' ;
  record(TID2, T2, T2', (version(X, VX, TSX, MDX), version(Y, VY', TSY', MDY')), WS2, true) ; LOG''
  = true if TSY' < TSY .
ceq fracRead(LOG ; record(TID2, ...) ; LOG' ; record(TID1, ...) ; LOG'') = true if TSY' < TSY .
eq fracRead(LOG) = false [owise] .

```

⁹ The additions to the original rule are written in italics.

The function `abortedRead` checks whether a transaction `TID2` reads a version `TSX` that was written by an aborted (flag `false`) transaction `TID1`:

```
op abortedRead : Record -> Bool .
eq abortedRead(LOG ;
  record(TID1, T1, T1', RS1, (version(X, VX, TSX, MDX), VS), false) ; LOG' ;
  record(TID2, T2, T2', (version(X, VX, TSX, MDX), VS), WS2, true) ; LOG'') = true .
eq abortedRead(LOG ; record(TID2, ...) ; LOG' ; record(TID1, ...) ; LOG'') = true.
eq abortedRead(LOG) = false [owise] .
```

No Lost Updates. We analyze the PLU property by searching for a final state in which the monitor shows that an update was lost:

```
search [1] initConfig =>! C:Config < M:Address : Monitor | log: LOG:Record >
  such that lu(LOG) .
```

The function `lu` checks whether there are lost updates in `LOG`. Lost updates concerns about conditional writes, namely, transactions fetching the same data. Once one of those transactions commits its writes, the others has to abort. Our specification of `lu` captures this by checking whether there are two transactions in *log* reading the same data (`sameReads`), and they both commit their writes on the same key(s) (`sameKeys`):

```
op lu : Record -> Bool .
ceq lu(LOG ; record(TID1, T1, T1', RS1, WS1, true) ; LOG' ;
  record(TID2, T2, T2', RS2, WS2, true) ; LOG'') = true
  if sameReads(RS1, RS2) /\ sameKeys(WS1, WS2) .
eq lu(LOG) = false [owise] .
```

We have performed our analysis with 4 different initial states, with up to 8 transactions, 2 data items and 4 partitions, without finding a violation of RA or PLU. We have also model checked the causal consistency (CC) property with the same initial states, and found a counterexample showing that ROLA does *not* satisfy CC. Each analysis command took about 30 seconds to execute on a 2.9 GHz Intel 4-Core i7-3520M CPU with 3.7 GB memory.

6 Statistical Model Checking of ROLA and Walter

The weakest consistency model in [7] guaranteeing RA and PLU is PSI, and the main system providing PSI is Walter [26]. ROLA must therefore outperform Walter to be an attractive design option. To quickly check whether ROLA does so, we have also modeled Walter—without its data replication features—in Maude (see <https://sites.google.com/site/fase18submission/maude-spec>), and use statistical model checking with PVESTA to compare the performance of ROLA and Walter in terms of throughput and average transaction latency.

Extracting Performance Measures from Executions. PVEStA estimates the expected (average) value of an expression on a run, up to a desired statistical confidence. The key to perform statistical model checking is therefore to define a measure on runs. Using the monitor in Section 5 we can define a number of functions on (states with) such a monitor that extract different performance metrics from this “system execution log”.

The function `throughput` computes the number of committed transactions per time unit. `committedNumber` computes the number of committed transactions in LOG, and `totalRunTime` returns the time when all transactions are finished (i.e., the largest *finishTime* in LOG):

```
op throughput : Config -> Float [frozen] .
eq throughput(< M : Monitor | log: LOG > REST)
  = committedNumber(LOG) / totalRunTime(LOG) .
```

The function `avgLatency` computes the average transaction latency by dividing the sum of the latencies of all committed transactions by the number of such transactions:

```
op avgLatency : Config -> Float [frozen] .
eq avgLatency(< M : Monitor | log: LOG > REST)
  = totalLatency(LOG) / committedNumber(LOG) .
```

where `totalLatency` uses the auxiliary function `$totalLatency` to compute the sum of all transaction latencies (time between the issue time and the finish time of a committed transaction).

```
op totalLatency : Record -> Float .
op $totalLatency : Record Float -> Float .
eq totalLatency(LOG) = $totalLatency(LOG,0.0) .
eq $totalLatency((record(TID,T1,T2,READS,WRITES,true) ; LOG),LATENCY) =
  $totalLatency(LOG,LATENCY + T2 - T1) .
eq $totalLatency((record(TID,T1,T2,READS,WRITES,false) ; LOG),LATENCY) =
  $totalLatency(LOG,LATENCY) .
eq $totalLatency(noRecord,LATENCY) = LATENCY .
```

Generating Initial States. We use an operator `init` to *probabilistically* generate initial states: `init(rtx, wtx, rwtx, part, keys, rops, wops, rwops, distr)` generates an initial state with *rtx* read-only transactions, *wtx* write-only transactions, *rwtx* read-write transactions, *part* partitions, *keys* data items, *rops* operations per read-only transaction, *wops* operations per write-only transaction, *rwops* operations per read-write transactions, and *distr* the key access distribution (the probability that an operation accesses a certain data item). To capture the fact that some data items may be accessed more frequently than others, we also use Zipfian distributions in our experiments.

Each PVEStA simulation starts from `init`, which rewrites to a *different* initial state in each simulation. The reason is that this expression involves generating certain values—such as the transactions—probabilistically.

`init` is defined by first generating the table, scheduler and monitor:

```

op init : NzNat NzNat NzNat NzNat NzNat
         NzNat NzNat NzNat KeyAccessDistr -> Config .
op $init : NzNat NzNat NzNat NzNat NzNat KeyVars KeyVars
         NzNat NzNat NzNat KeyAccessDistr Config -> Config .
op $$init : NzNat NzNat NzNat NzNat NzNat KeyVars KeyVars
         NzNat NzNat NzNat KeyAccessDistr Nat Config -> Config .

eq init(RTX,WTX,RWTX,PS,KS,ROPS,WOPS,RWOPS,KAD) =
  $init(RTX,WTX,RWTX,PS,PS,kvars(KS,keyVars),kvars(KS,keyVars),
        ROPS,WOPS,RWOPS,KAD, < 0 . 1 : Table | table: [emptyTable] >)
  { 0.0 | nil } < 0 . 2 : Monitor | log: noRecord > .

```

where $\$init$ and $$$init$ are two auxiliary functions which continue to generate and update other objects. $kvars$ cuts out the first KS number of key-local var pairs, $\langle k_1, k_1l \rangle ; \langle k_2, k_2l \rangle ; \dots ; \langle k_{ks}, k_{ks}l \rangle$, from all constant key-local pairs.

Then $\$init$ uniformly assigns each key to a partition; $assign$ also updates the table with the assigned key and its partition:

```

eq $init(RTX,WTX,RWTX,0,REPLS2,< K,VAR > ; KVARs),KVARs',
        ROPS,WOPS,RWOPS,KAD,C) =
  $init(RTX,WTX,RWTX,0,REPLS2,KVARs,KVARs',ROPS,WOPS,RWOPS,KAD,
        assignKey(K,sampleUniWithInt(REPLS2) + 1,C)) .

op assignKey : Key Address Config -> Config .
eq assignKey(K,PID,< PID : Partition | datastore: VS, AS >
  < TB : Table | table: [KEYREPLICAS] > C) =
  < PID : Partition | datastore: (VS version(K,[0],eptTS,empty)), AS >
  < TB : Table | table: [replicatingSites(K,PID) ;; KEYREPLICAS] > C .

```

where $eptTS$ is the default timestamp.

Then $\$init$ continues to generate transactions when all keys have been assigned (denoted by $noKeyVar$). The following case shows when the remaining transactions have all three types (denote by $s\ RTX$, $s\ WTX$ and $s\ RWTX$):

```

eq $init(s RTX,s WTX,s RWTX,0,PS',noKeyVar,KVARs',
        ROPS,WOPS,RWOPS,KAD,C) =
  $$init(s RTX,s WTX,s RWTX,0,PS',noKeyVar,KVARs',ROPS,WOPS,RWOPS,
        KAD,sampleUniWithInt(s RTX + s WTX + s RWTX),C) .

eq $$init(s RTX,s WTX,s RWTX,0,PS',noKeyVar,KVARs',ROPS,WOPS,RWOPS,
        KAD,R-OR-W-OR-RW,C) =
  if R-OR-W-OR-RW < s RTX --- generate read-only txns
  then $init(RTX,s WTX,s RWTX,0,PS',noKeyVar,KVARs',
        ROPS,WOPS,RWOPS,KAD,addRTxn(sampleUniWithInt(PS') + 1,
        ROPS,KVARs',KAD,C))

```

```

else if s RTX <= R-OR-W-OR-RW and R-OR-W-OR-RW < s RTX + s WTX
  --- generate write-only txns
  then $init(s RTX, WTX, s RWTX, 0, PS', noKeyVar, KVARs', ROPS,
    WOPS, RWOPS, KAD, addWTxn(sampleUniWithInt(PS') + 1,
    WOPS, KVARs', KAD, C))
  else --- generate read-write txns
  $init(s RTX, s WTX, RWTX, 0, PS', noKeyVar, KVARs', ROPS,
    WOPS, RWOPS, KAD, addRWTxn(sampleUniWithInt(PS') + 1,
    RWOPS, KVARs', KAD, C))
fi fi .

```

This rule first probabilistically decides whether the next transaction is a read-only, a write-only, or a read-write transaction. Since the probability of picking a read transaction should be $\frac{\#readsLeft}{\#txnLeft}$, it uniformly picks a value R-OR-W-OR-RW from $[0, \dots, \#txnLeft - 1]$ (the number of transactions left to generate is $s RTX + s WTX + s RWTX$) using the expression `sampleUniWithInt(s RTX + s WTX + s RWTX)`. If the value picked is in $[0, \dots, \#readsLeft - 1]$ ($< s RTX$), we generate a new read-only transaction next (`then` branch); otherwise, in a similar way, we generate a new write-only transaction (`else-if-then` branch), or a new read-write transaction (`else-if-else` branch). But which partition should get the transaction? The partitions have identities 1, 2, ..., n, where n is the number of partitions (PS'). The expression `sampleUniWithInt(PS') + 1` gives us the partition, sampled uniformly from $[1, \dots, n]$.

Similarly, we treat other cases based on the type(s) of the remaining transactions. We omit the details for simplicity.

The following defines the functions `addReadTrans` which generates a new read-only transaction:

```

op addRTxn : Address Nat KeyVars KeyAccessDistr Config -> Config .
op $addRTxn : Config -> Config . --- generate local variables

--- if this is the first ro-txn to generate
eq addRTxn(RID, ROPS, KVARs, KAD, < PID : Partition | gotTxns: emptyTxnList, AS > C) =
  $addRTxn(< PID : Partition | gotTxns: < PID . 1 : Txn |
    operations: addReads(ROPS, KVARs, KAD), readSet: nil, latest: empty,
    localVars: empty, txnSQN: 0 >, AS >) C .

--- if there is already some txn(s) generated
eq addRTxn(RID, ROPS, KVARs, KAD, < PID : Partition | gotTxns:
  (TXNS ;; < PID . N : Txn | AS' >), AS > C) =
  $addRTxn(< PID : Partition | gotTxns: (TXNS ;; < PID . N : Txn | AS' >
    ;; < PID . (N + 1) : Txn | operations: addReads(ROPS, KVARs, KAD), readSet: nil,
    latest: empty, localVars: empty, txnSQN: 0 >), AS >) C .

--- pair local variables
eq $addRTxn(< PID : Partition | gotTxns: (TXNS ;; < PID . N : Txn |
  operations: OPS, readSet: nil, latest: empty, localVars: empty,
  txnSQN: 0 >), AS >) =

```



```
< PID : Partition | gotTxns: (TXNS ;; < PID . N : Txn | operations: OPS,
readSet: nil, latest: empty, localVar: lvars(OPS), txnSQN: 0 >), AS > .
```

where `lvars` generates the local variables by projecting the associated local variable from each key-local variable pair. We omit the cases of `addWTxn` and `addRWTxn` for simplicity.

When there are no more transactions to generate, `$init` returns the generated objects:

```
eq $init(0,0,0,0,PS',noKeyVar,KVARS',ROPS,WOPS,RWOPS,KAD,C) = C .
```

Statistical Model Checking Results. We performed our experiments under different configurations, with 200 transactions, 2–4 operations per transaction, up to 200 data items and up to 50 partitions, with lognormal message delay distributions, and with uniform and Zipfian data item access distributions. Regarding Lognormal’s parameters, local delays use $\mu = 0$ and $\sigma = 1$, while remote delays use $\mu = 3$ and $\sigma = 2$.

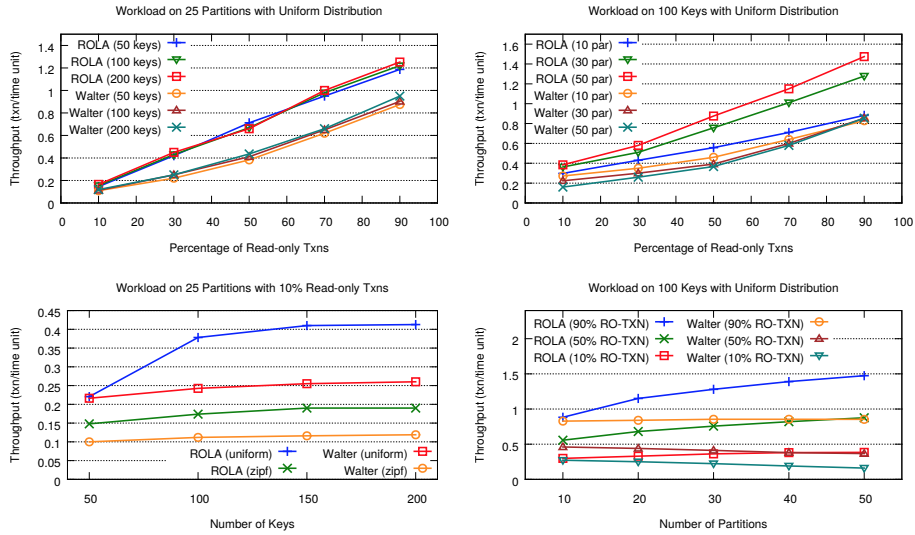


Fig. 2. Throughput comparison under different workload conditions.

The plots in Fig. 2 show the *throughput* as a function of the percentage of read-only transactions, number of partitions, and number of keys (data items), sometimes with both uniform and Zipfian distributions. The plots show that ROLA outperforms Walter for all parameter combinations. More partitions gives ROLA higher throughput (since concurrency increases), as opposed to Walter

(since Walter has to propagate transactions to more partitions to advance the vector timestamp). We only plot the results under uniform key access distribution, which are consistent with the results using Zipfian distributions.

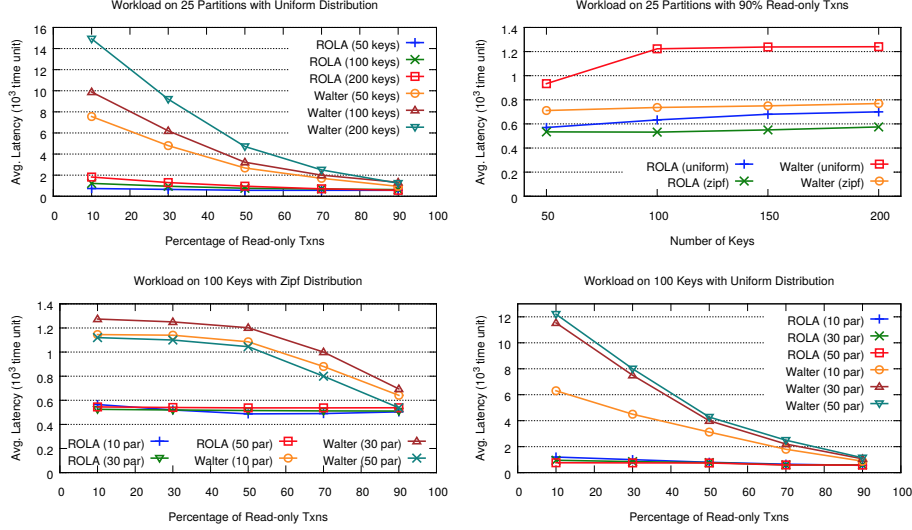


Fig. 3. Average latency comparison across varying workload conditions.

The plots in Fig. 3 show the *average transaction latency* as a function of the same parameters as the plots for throughput. Again, we see that ROLA outperforms Walter in all settings. In particular, this difference is quite large for write-heavy workloads; the reason is that Walter incurs more and more overhead for providing causality, which requires background propagation to advance the vector timestamp. The latency tends to converge under read-heavy workload (because reads in both ROLA and Walter can commit locally without certification), but ROLA still has noticeable lower latency than Walter.

ROLA is approximately 850 LOC, while Walter 1200 LOC (all excluding approximate 300 shared LOC related to the scheduler and sampler, and 350 shared LOC related to the initial-states generator). Computing the probabilities took 6 hours (worst case) on 10 servers, each with a 64-bit Intel Quad Core Xeon E5530 CPU with 12 GB memory. Each point in the plots represents the average of three statistical model checking results.

7 Related Work

Maude and PVESTA have been used to model and analyze the correctness and performance of a number of distributed data stores: the Cassandra key-value

store [21,17,18], different versions of RAMP [19,16], the Walter transactional data store [20], and Google’s Megastore [10,11]. In contrast to these papers, our paper uses formal methods to develop and validate an entirely new design, ROLA, for a new consistency model.

We are not aware of other work on formal model-based performance analysis of globally-distributed transactional databases. This might be because the most popular formal tools supporting probabilistic/statistical model checking are based on automata (e.g., Uppaal SMC [2] and Prism [1]), and it seems hard, or impossible, to model state-of-the-art distributed transactional systems using automata. Maude provides the expressiveness and modeling convenience that makes it possible to model such complex systems with reasonable effort.

Concerning formal methods for distributed data stores, engineers at Amazon have used TLA+ and its model checker TLC to model and analyze the correctness of key parts of Amazon’s celebrated cloud computing infrastructure [23]. In contrast to our work, they only use formal methods for correctness analysis; indeed, one of their complaints is that they cannot use their formal method for performance estimation. The designers of the TAPIR transaction protocol for distributed storage systems have also specified and model checked correctness (but not performance) properties of their design using TLA+ [29].

In contrast to our work that aims at developing and analyzing high-level formal models to quickly analyze different design choices and finding bugs early, other approaches [15,27] use distributed model checkers to analyze the *implementation* of cloud storage systems. Verifying both protocols and code is the goal of the IronFleet framework at Microsoft Research [12]. Their verification methodology includes a wide range of methods and tools, and requires (in contrast to our method) “considerable assistance from the developer.”

8 Conclusions

We have presented the formal design and analysis of ROLA, a distributed transaction protocol that supports a new consistency model not present in the survey by Cerone *et al.* [7]. Using formal modeling and both standard and statistical model checking analyses we have: (i) validated ROLA’s RA and PLU consistency requirements; and (ii) analyzed its performance requirements, showing that ROLA outperforms Walter in all performance measures.

This work has shown, to the best of our knowledge for the first time, that the design and validation of a *new* distributed transaction protocol can be achieved relatively quickly *before* its implementation by the use of formal methods. Our next planned step is to implement ROLA, evaluate it experimentally, and compare the experimental results with the formal analysis ones. In previous work on existing systems such as Cassandra [13], RAMP [5] and Walter [26], the performance estimates obtained by formal analysis and those obtained by experimenting with the real system were basically in agreement with each other [17,16,20]. This confirmed the useful predictive power of the formal analyses. Our future research will investigate the existence of a similar agreement for ROLA. This

work is part of a long-term research effort in which we have been using Maude to both meet the challenges and exploit the opportunities of modular design and analysis for cloud-based transaction systems (see [6] for a survey).

References

1. PRISM, <http://www.prismmodelchecker.org/>
2. Uppaal SMC, <http://people.cs.aau.dk/~adavid/smc/>
3. Agha, G.A., Meseguer, J., Sen, K.: PMAude: Rewrite-based specification language for probabilistic object systems. *Electr. Notes Theor. Comput. Sci.* 153(2) (2006)
4. Alturki, M., Meseguer, J.: PVeStA: A parallel statistical model checking and quantitative analysis tool. In: CALCO’11. LNCS, vol. 6859. Springer (2011)
5. Bailis, P., Fekete, A., Ghodsi, A., Hellerstein, J.M., Stoica, I.: Scalable atomic visibility with RAMP transactions. *ACM Trans. Database Syst.* 41(3), 15:1–15:45 (2016)
6. Bobba, R., Grov, J., Gupta, I., Liu, S., Meseguer, J., Ölveczky, P.C., Skeirik, S.: Design, formal modeling, and validation of cloud storage systems using Maude. Tech. rep., University of Illinois at Urbana-Champaign (2017), <http://hdl.handle.net/2142/96274>
7. Cerone, A., Bernardi, G., Gotsman, A.: A framework for transactional consistency models with atomic visibility. In: CONCUR. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik (2015)
8. Clavel, M., Durán, F., Eker, S., Lincoln, P., Martí-Oliet, N., Meseguer, J., Talcott, C.L.: All About Maude, LNCS, vol. 4350. Springer (2007)
9. Eckhardt, J., Mühlbauer, T., Meseguer, J., Wirsing, M.: Statistical model checking for composite actor systems. In: WADT’12. LNCS, vol. 7841. Springer (2013)
10. Grov, J., Ölveczky, P.C.: Formal modeling and analysis of Google’s Megastore in Real-Time Maude. In: Specification, Algebra, and Software. LNCS, vol. 8373. Springer (2014)
11. Grov, J., Ölveczky, P.C.: Increasing consistency in multi-site data stores: Megastore-CGC and its formal analysis. In: SEFM. LNCS, vol. 8702. Springer (2014)
12. Hawblitzel, C., Howell, J., Kapritsos, M., Lorch, J.R., Parno, B., Roberts, M.L., Setty, S.T.V., Zill, B.: IronFleet: proving practical distributed systems correct. In: Proc. 25th Symposium on Operating Systems Principles (SOSP’15). ACM (2015)
13. Hewitt, E.: Cassandra: The Definitive Guide. O’Reilly Media (2010)
14. Kallman, R., Kimura, H., Natkins, J., Pavlo, A., Rasin, A., Zdonik, S.B., Jones, E.P.C., Madden, S., Stonebraker, M., Zhang, Y., Hugg, J., Abadi, D.J.: H-store: a high-performance, distributed main memory transaction processing system. *PVLDB* 1(2), 1496–1499 (2008)
15. Leesatapornwongsa, T., Hao, M., Joshi, P., Lukman, J.F., Gunawi, H.S.: SAMC: Semantic-aware model checking for fast discovery of deep bugs in cloud systems. In: 11th USENIX Symposium on Operating Systems Design and Implementation (OSDI’14). USENIX Association (2014)
16. Liu, S., Ölveczky, P.C., Ganhotra, J., Gupta, I., Meseguer, J.: Exploring design alternatives for RAMP transactions through statistical model checking. In: Proc. ICFEM’17. LNCS, vol. 10610. Springer (2017)
17. Liu, S., Ganhotra, J., Rahman, M., Nguyen, S., Gupta, I., Meseguer, J.: Quantitative analysis of consistency in NoSQL key-value stores. *Leibniz Transactions on Embedded Systems* 4(1), 03:1–03:26 (2017)

18. Liu, S., Nguyen, S., Ganhotra, J., Rahman, M.R., Gupta, I., Meseguer, J.: Quantitative analysis of consistency in NoSQL key-value stores. In: QEST. pp. 228–243 (2015)
19. Liu, S., Ölveczky, P.C., Rahman, M.R., Ganhotra, J., Gupta, I., Meseguer, J.: Formal modeling and analysis of RAMP transaction systems. In: SAC’16. ACM (2016)
20. Liu, S., Ölveczky, P.C., Santhanam, K., Wang, Q., Gupta, I., Meseguer, J.: Formal modeling and analysis of the Walter transactional data store. Tech. rep., <http://hdl.handle.net/2142/98988>
21. Liu, S., Rahman, M.R., Skeirik, S., Gupta, I., Meseguer, J.: Formal modeling and analysis of Cassandra in Maude. In: ICFEM’14. LNCS, vol. 8829. Springer (2014)
22. Meseguer, J.: Conditional rewriting logic as a unified model of concurrency. *Theoretical Computer Science* 96(1), 73–155 (1992)
23. Newcombe, C., Rath, T., Zhang, F., Munteanu, B., Brooker, M., Deardeuff, M.: How Amazon Web Services uses formal methods. *Communications of the ACM* 58(4), 66–73 (2015)
24. Sen, K., Viswanathan, M., Agha, G.: On statistical model checking of stochastic systems. In: CAV’05. LNCS, vol. 3576. Springer (2005)
25. Sen, K., Viswanathan, M., Agha, G.A.: VESTA: A statistical model-checker and analyzer for probabilistic systems. In: QEST’05. IEEE Computer Society (2005)
26. Sovran, Y., Power, R., Aguilera, M.K., Li, J.: Transactional storage for geo-replicated systems. In: SOSP 2011. ACM (2011)
27. Yang, J., Chen, T., Wu, M., Xu, Z., Liu, X., Lin, H., Yang, M., Long, F., Zhang, L., Zhou, L.: MODIST: transparent model checking of unmodified distributed systems. In: Proc. 6th USENIX Symposium on Networked Systems Design and Implementation (NSDI’09). pp. 213–228. USENIX Association (2009)
28. Younes, H.L.S., Simmons, R.G.: Statistical probabilistic model checking with a focus on time-bounded properties. *Inf. Comput.* 204(9), 1368–1409 (2006)
29. Zhang, I., Sharma, N.K., Szekeres, A., Krishnamurthy, A., Ports, D.R.K.: Building consistent transactions with inconsistent replication. In: Proc. Symposium on Operating Systems Principles, (SOSP’15). ACM (2015)

A The RAMP-Fast Algorithm as Given in [5]

Figure 4 shows the RAMP-Fast algorithm as it is described in [5].

ALGORITHM 1: RAMP-Fast**Server-side Data Structures**

- 1: *versions*: set of versions (*item*, *value*, timestamp ts_v , metadata *md*)
- 2: *lastCommit*[*i*]: last committed timestamp for item *i*

Server-side Methods

- 3: **procedure** PREPARE(*v* : version)
- 4: *versions*.add(*v*)
- 5: **return**
- 6: **procedure** COMMIT(ts_c : timestamp)
- 7: $I_{ts} \leftarrow \{w.item \mid w \in versions \wedge w.ts_v = ts_c\}$
- 8: $\forall i \in I_{ts}, lastCommit[i] \leftarrow \max(lastCommit[i], ts_c)$
- 9: **procedure** GET(*i* : item, ts_{req} : timestamp)
- 10: **if** $ts_{req} = \emptyset$ **then**
- 11: **return** $v \in versions : v.item = i \wedge v.ts_v = lastCommit[item]$
- 12: **else**
- 13: **return** $v \in versions : v.item = i \wedge v.ts_v = ts_{req}$

Client-side Methods

- 14: **procedure** PUT_ALL(*W* : set of (*item*, *value*))
- 15: $ts_{tx} \leftarrow$ generate new timestamp
- 16: $I_{tx} \leftarrow$ set of items in *W*
- 17: **parallel-for** (*i*, *v*) $\in W$
- 18: $w \leftarrow \langle item = i, value = v, ts_v = ts_{tx}, md = (I_{tx} - \{i\}) \rangle$
- 19: invoke PREPARE(*w*) on respective server (i.e., partition)
- 20: **parallel-for** server *s* : *s* contains an item in *W*
- 21: invoke COMMIT(ts_{tx}) on *s*
- 22: **procedure** GET_ALL(*I* : set of items)
- 23: $ret \leftarrow \{\}$
- 24: **parallel-for** *i* $\in I$
- 25: $ret[i] \leftarrow GET(i, \emptyset)$
- 26: $v_{latest} \leftarrow \{\}$ (default value: -1)
- 27: **for** response *r* $\in ret$ **do**
- 28: **for** $i_{tx} \in r.md$ **do**
- 29: $v_{latest}[i_{tx}] \leftarrow \max(v_{latest}[i_{tx}], r.ts_v)$
- 30: **parallel-for** item *i* $\in I$
- 31: **if** $v_{latest}[i] > ret[i].ts_v$ **then**
- 32: $ret[i] \leftarrow GET(i, v_{latest}[i])$
- 33: **return** *ret*

Fig. 4. The RAMP-Fast algorithm as described in [5].